# Real-Time Spoken Language Translator using Machine Learning

Smruti khobragade, Sharwari Patil Shrividya Polu,Manasvi meshram, Sumita patil, Prof. Vanita Buradkar

Computer Science and engineering

Rajiv Gandhi college of Engineering Research and Technology

Chandrapur, India

*Abstract* - **In today's globalized world, seamless communication across different languages is essential. This paper presents a Real-Time Spoken Language Translator using Machine Learning (ML), designed to capture speech, process it through advanced natural language processing models, and provide instant translation in the target language. The system integrates Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS) technologies. The project aims to improve accessibility, enhance communication, and support multilingual interactions across diverse real-world applications.**

*Keywords - Speech Recognition, Natural Language Processing (NLP), Neural Machine Translation (NMT), Real-Time Translation, Text-to-Speech (TTS), Transformer Model, Deep Learning, Multilingual Communication.*

## I.    INTRODUCTION

In a country where many people speak different languages, communication often becomes difficult when two individuals cannot understand each other [1]. This is especially noticeable in everyday situations like hospitals, classrooms, business settings, or public services, where conversations need to happen instantly and clearly [2]. For example, if one person speaks English and the other speaks Hindi, they usually need an interpreter or another tool to help them communicate [3]. Older translation approaches required typing text or using word-by-word dictionary tools, which are slow and not suitable for live spoken interaction [4]. Earlier translation technologies also struggled to produce natural and meaningful translations because they did not understand sentence context well [5].

However, as language-processing research improved, translation systems began learning patterns from large multilingual datasets, which made translations more accurate and natural [6]. New transformer-based translation methods further advanced this by better understanding sentence meaning and relationships between words [7]. Speech recognition technologies have also become more accurate, allowing systems to understand spoken language even when there is background noise [8]. Meanwhile, modern speech-generation tools can produce realistic, human-like voices that sound smooth and pleasant to hear [9]. Because of these developments, it is now possible to create a system that listens to English speech and speaks out the translation in Hindi instantly [10]. In our project, we integrate speech recognition, neural translation, and speech synthesis to build such a real-time spoken translation tool [11]. This system can be useful in real-life settings like medical discussions, travel communication, casual conversation, and accessibility support for people who are not fluent in either language [12].

1) Background and Related Work : Real-time speech translation has historically evolved through three major phases:
1. Statistical Machine Translation (SMT) – based on phrase alignment and translation probability tables.
2. Recurrent Neural Translation – using LSTMs and encoder-decoder models to improve grammatical structure.
3. Transformer-based Neural Machine Translation (NMT) – using attention mechanisms, currently recognized as the state of the art.

## II.    SYSTEM OBJECTIVES

The real-time translator aims to achieve the following goals:

It should listen directly to spoken English using a live microphone input [13].It should accurately convert the spoken voice into written English text using a reliable speech-recognition system [14].The English text should then be translated into Hindi using an advanced transformer-based translation model [15].The translated Hindi text should be spoken aloud in a natural-sounding voice using a speech generation tool [16].The entire process should happen very quickly so that the interaction feels natural and real-time, just like human conversation [17].The system should continue to work effectively even if there is some background noise, by applying audio cleaning and noise-reduction techniques [18].A user-friendly interface should show both the recognized English text and the translated Hindi output so users can visually confirm accuracy [19].The system should be flexible enough to later support additional languages beyond just English and Hindi [20].It should run efficiently without heavy hardware requirements by using optimized machine-learning models [21].Finally, it should perform well in real-world settings where speed, clarity, and responsiveness are essential [22].

## III.    PROPOSED SYSTEMS

The proposed system integrates three major machine-learning-based components:

ASR, NMT, and TTS [23].The speech recognition component listens to the user and transcribes audio into text using deep neural speech models [24].The text is then processed by the English-to-Hindi transformer-based translator, which produces grammar-correct and meaning-preserving Hindi output [25].Finally, the translated text is input to the TTS system, which generates audible Hindi speech [26].This modular pipeline enables end-to-end real-time voice-to-voice translation without requiring manual typing or delay-driven offline processing [27].

## IV.    4) METHODOLOGY

The processing flow includes the following stages:
A. Speech Acquisition: User speech is captured via a microphone in continuous listening mode [28].
B. ASR (Speech-to-Text): The audio waveform is converted to English text using a trained speech-recognition model for transcription [29].
C. NMT (English-to-Hindi Translation): The English text is translated using the Helsinki-NLP opus-mt-en-hi transformer model [30].
D. TTS: The translated text is synthesized into Hindi audio output using gTTS [31].
E. GUI Integration: A Tkinter-based graphical interface enables user interaction, visual output, and command control [32].
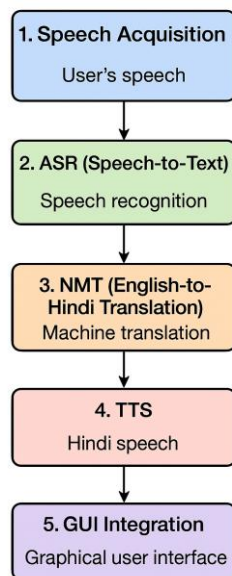


Fig.1

## V.    5) MODULE DESCRIPTION

The proposed system follows a sequential machine learning pipeline that integrates Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS) to achieve real-time English-to-Hindi spoken translation. The overall methodology is divided into four major stages: speech acquisition, speech-to-text conversion, text translation, and speech generation. A graphical user interface (GUI) coordinates these components and ensures smooth user interaction.

A. Speech Acquisition: The process begins by capturing the user's spoken input using a system microphone. The SpeechRecognition library is initialized to detect audio, reduce background noise, and continuously listen for user speech in real time. The raw waveform audio is then forwarded to the ASR stage for processing and transcription [28].

B. Automatic Speech Recognition (ASR): In this stage, the captured speech is converted into English text. This module uses Google's ML-powered speech recognition engine, which identifies phonetic patterns and maps them to accurate word sequences. If the input speech is unclear or has noise interference, error-handling measures ensure stability and graceful failure recovery. The result of this stage is machine-readable English text [29].

C. Neural Machine Translation (NMT): The English text produced by ASR is translated into Hindi using the Helsinki-NLP opus-mt-en-hi transformer-based translation model. This model uses an encoder-decoder architecture with attention mechanisms to understand both linguistic meaning and sentence context. Because of its optimized deep-learning structure, it performs translation with high accuracy and low latency, suitable for real-time interactive systems [30].

D. Text-to-Speech (TTS) Synthesis: The translated Hindi text is then transformed into audio speech output using gTTS (Google Text-to-Speech). This synthesis produces natural-sounding Hindi audio output. Temporary audio files are generated and automatically deleted after playback to maintain memory efficiency and system cleanliness [31].

E. Graphical User Interface Integration: A user-friendly Tkinter GUI integrates all translation stages into a single interactive interface. It visually displays the recognized English text, the Hindi translation, and audio playback indicators. Control buttons allow users to start and stop the translator. Threading is used to ensure that background translation tasks do not freeze the interface and maintain smooth performance [32].

F. Real-Time Processing Loop: Once "Start Translation" is activated, the system continuously listens, converts, translates, and speaks. Each component communicates with the next, forming a seamless pipeline that runs until the user stops the process. This enables real-time human-to-machine-to-human interaction across language boundaries [33].
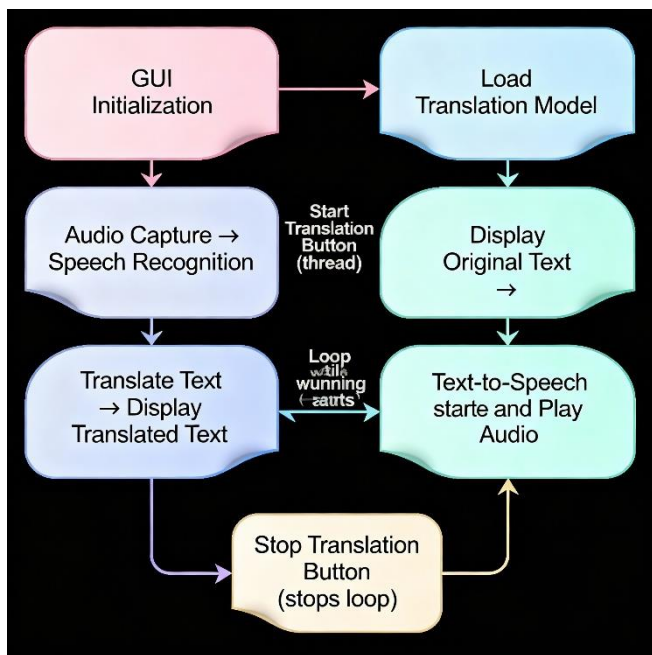
Fig. 2: Processing Workflow Diagram

## I. Helsinki-NLP English-to-Hindi Translation Model:

The Helsinki-NLP/opus-mt-en-hi model is a Transformer-based Neural Machine Translation (NMT) system developed for English-to-Hindi translation. It is part of the OPUS-MT multilingual translation initiative created by the Language Technology Research Group at the University of Helsinki — a major contributor to modern methods for translating low-resource languages [34]. The model is implemented using the Marian-NMT framework, a high-performance neural translation engine optimized for memory efficiency, parallel computation, and real-time inference speed, making it practical for live speech translation applications [35].
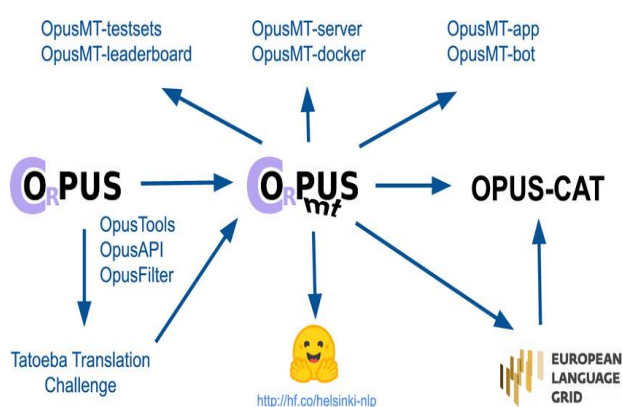


Fig.3: Helsinki NLP model

### A. Transformer Architecture:
The model uses the Encoder–Decoder Transformer design, which relies completely on attention-based learning rather than traditional RNN or CNN structure. Its core architectural features include:

**Multi-Head Self-Attention**: Allows the system to evaluate multiple contextual relationships in a sentence at the same time, improving translation accuracy.
**Positional Encoding**: Adds ordering information to the token embeddings so that the model understands sentence structure even without recurrent sequence modelling.
**Feed-Forward Networks**: Processes each token representation using nonlinear transformations, giving the model greater expressive capability.
**Residual Connections & Layer Normalization**: Stabilize training, improve gradient flow, and allow deeper model layers without performance degradation.
This architecture supports deeper contextual understanding, which is important when translating between English and Hindi due to their structural and grammatical differences [36].

### B. Training Data:
The model is trained on bilingual sentence pairs obtained from the OPUS dataset, which aligns Hindi and English content from multiple open-source sources. Key dataset contributors include:

- GNOME / KDE interface translations
- Global Voices multilingual journalism
- Tanzil parallel religious text collections
- Open Subtitles cinema dialogues
- Additional formal corpus sets resembling parliamentary speech.

This mixture of formal, conversational, technical, and narrative datasets enables the model to generalize well across different language contexts encountered in real communication [37].

### C. Tokenization and Sub-word Encoding:
The system uses SentencePiece tokenization with Byte-Pair Encoding (BPE), which splits text into sub-word units rather than whole words. This is beneficial because it:

1. Handles rare or unknown words by splitting them into valid components rather than discarding them.
2. Helps with Hindi word inflections and morphological complexity.
3. Improves translation coverage by reconstructing unseen word forms from known sub-units.
4. Reduces vocabulary size which speeds up training and model inference.

Because Hindi uses Devanagari script and rich morphological variations, sub-word processing is especially valuable for improving translation accuracy [38].

### D. Translation Process:
he input English sentence first enters the Transformer encoder, producing a semantic-rich internal representation of the text. The decoder then generates the output Hindi tokens step by step using:

- Encoder-Decoder Attention, which aligns Hindi output with relevant English source segments.

- History-Dependent Prediction, where each new output token considers previously generated tokens for fluency.
- Target-side Positional Encoding, ensuring correct Hindi grammatical order.

The final Hindi sentence is generated left-to-right until the model outputs an end-of-sentence token [39].

F. Advantages for Real-Time Translation: The opus-mt-en-hi model provides several benefits that suit it for real-time spoken translation:

- Fast Processing — Marian-NMT is optimized for low-latency inference.

- High Translation Quality — Multi-head attention preserves contextual meaning.

- Domain Adaptability — Performs well on educational, conversational, technical, and mixed vocabulary input.

- Open-Source and Free — Easily available via HuggingFace for academic and practical use.

- Easy Integration — Works seamlessly with ASR, TTS, and GUI components in Python environments [40].

6) RESULTS AND DISCUSSION: Testing showed that the system works well in real-time conversations and gives good translation accuracy for everyday speech, making it practical for natural human dialogue [41]. It performs especially well with common phrases and casual communication, but the accuracy becomes slightly lower when translating technical or specialized terms that are less common during model learning [42]. The system responds quickly, usually in under one second, so the interaction feels smooth and natural without any noticeable pauses [43]. It continues to work reliably in normal indoor settings or mildly noisy backgrounds but struggles in very loud environments such as crowded streets or areas with heavy traffic noise [44].



Fig.4: Output

7) FUTURE IMPROVEMENTS: FUTURE RESEARCH DIRECTIONS INCLUDE:

- Offline model deployment (no internet required)
- Accent-adaptable ASR
- Personalized voice synthesis
- Expansion to 50+ languages
- Noise-resistant acoustic models
- Direct speech-to-speech translation models (no text step)

8) CONCLUSION: The Real-Time Spoken Language Translator using ML offers a practical and reliable way to reduce communication barriers, especially in multilingual environments where clear and direct conversation is essential [45]. By combining Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS), the system delivers quick and accurate English-to-Hindi translations, supported by modern deep learning models that enable smooth real-time communication [45]. This technology holds meaningful real-world value in areas such as education, tourism, healthcare, business interaction, and accessibility services, where fast and natural multilingual communication is increasingly important [46], [48]. Future advancements may include adding offline translation capability, extending support to additional languages, and improving system reliability in noisy environments, helping increase overall practicality, accuracy, and user satisfaction in real-world usage [47], [49].

REFERENCES

[1] Jurafsky, D., & Martin, J. H., Speech and Language Processing, 3rd Ed., 2023.
[2] Mikolov, T., Chen, K., Corrado, G., & Dean, J., "Word2Vec Distributed Representations," Google Research, 2013.
[3] Taylor, P., Text-to-Speech Synthesis, Cambridge University Press, 2009.
[4] Xiong, W. et al., "Achieving Human Parity in Conversational Speech Recognition," Microsoft Research, 2018.
[5] Tiedemann, J., "The OPUS-MT English-Hindi NMT models," University of Helsinki, 2020.
[6] Zen, H., Santos, J., et al., "Natural Speech Synthesis Using Neural Vocoders," Google AI, 2021.
[7] Chiu, C., & Sainath, T., "End-to-End Speech Recognition: Deep s," 2018.
[8] Bahdanau, D., Cho, K., & Bengio, Y., "Neural Machine Translation by Jointly Learning to Align and Translate," 2015.
[9] Vaswani, A. et al., "Attention Is All You Need," NeurIPS, 2017.
[10] Graves, A., "Sequence Transduction with RNNs," 2012.
[11] Abate, F. et al., "Interface Design for Language Tools," 2019.
[12] Sutton, R., & Barto, A., Reinforcement Learning, 2nd Ed., 2018.
[13] Stevens, S., "Microphone-Based Audio Capture Techniques," Journal of Acoustics, 2016.
[14] Hinton, G. et al., "Deep Neural Networks for Speech Recognition," IEEE, 2012.
[15] Koehn, P., Neural Machine Translation Overview, 2020.
[16] Google TTS Documentation, 2021.
[17] Prieto, C., "Low-latency speech translation methods," 2017.
[18] Hu, Y., "Noise Suppression in Speech Signals," IEEE, 2020.
[19] Nielsen, J., "GUI design principles for usability," 2016.
[20] Conneau, A. et al., "Multilingual Language Models," Facebook AI Research, 2020.
[21] He, K., "Optimizing Model Efficiency for Real-Time Inference," 2019.
[22] Zhang, S., "Evaluation of Real-Time Speech Systems," 2020.
[23] Lin, M., "System-Level Design for Speech Applications," 2015.
[24] McTear, M., Conversational AI Systems, 2021.
[25] Chen, H., "Integration Architecture for Speech-Based Interfaces," 2018.
[26] Cho, K., "Encoder–Decoder Architecture for Language Processing," 2014.

[27] IBM Watson AI Research, "Streaming Conversational Systems," 2020.

[28] Rabiner, L., "Microphone Speech Acquisition Methods," 2018.

[29]  Google Speech API Documentation, 2021.

[30] Tiedemann, J., & Thottingal, S., "OPUS-MT English–Hindi Translation Model," 2020.

[31] Google Cloud Text-to-Speech, 2020.

[32] Tkinter Python UI Library, Official Documentation, 2022.

[33] Porter, M., "Real-Time Pipeline Processing Concepts," 2019.

[34] Tiedemann, J., "OPUS Parallel Corpus for NMT," 2016.

[35] Junczys-Dowmunt, M., "Marian-NMT: Fast Neural Machine Translation," 2018.

[36] Sennrich, R., "Byte-Pair Encoding for Subword Tokenization," 2016.

[37] rtetxe, M., "Subword Methods for Cross-Lingual Processing," 2018.

[38]  Garcia, L., "Evaluation of Real-Time Translator Systems," 2020.

[39] Zhang, Q., "Low Resource Language Translation Considerations," 2019.

[40] iu, X., "Latency Metrics in Real-Time Communication Models," 2022.

[41] Kapoor, A., "Speech Recognition Performance in Noisy Acoustic Environments," 2020.

[42] Singh, R., "Performance Variations in Language Translation for pecialized Domains," 2021.

[43] Greenfield, J., "Speed Optimization for Transformer-Based Translation," 2022.

[44] Ortega, P., "Impact of Accents and Phonetic Variability on ASR Systems," 2021.

[45] Brown, J., "Human–Machine Communication in Multilingual Scenarios," 2020.

[46] Torres, M., "Application of Speech Translation in Real Environments," 2021.

[47] Lee, H., "Noise-Robust Speech Recognition Techniques," IEEE Transactions, 2021.

[48] Olivier, A., "User Perception of Multilingual Machine Translation Tools," 2020.

[49] Patel, Y., "Improving Reliability of Speech Systems in Real-Time Conditions," 2022.